# The IBM High Performance Computing Toolkit on BlueGene/L

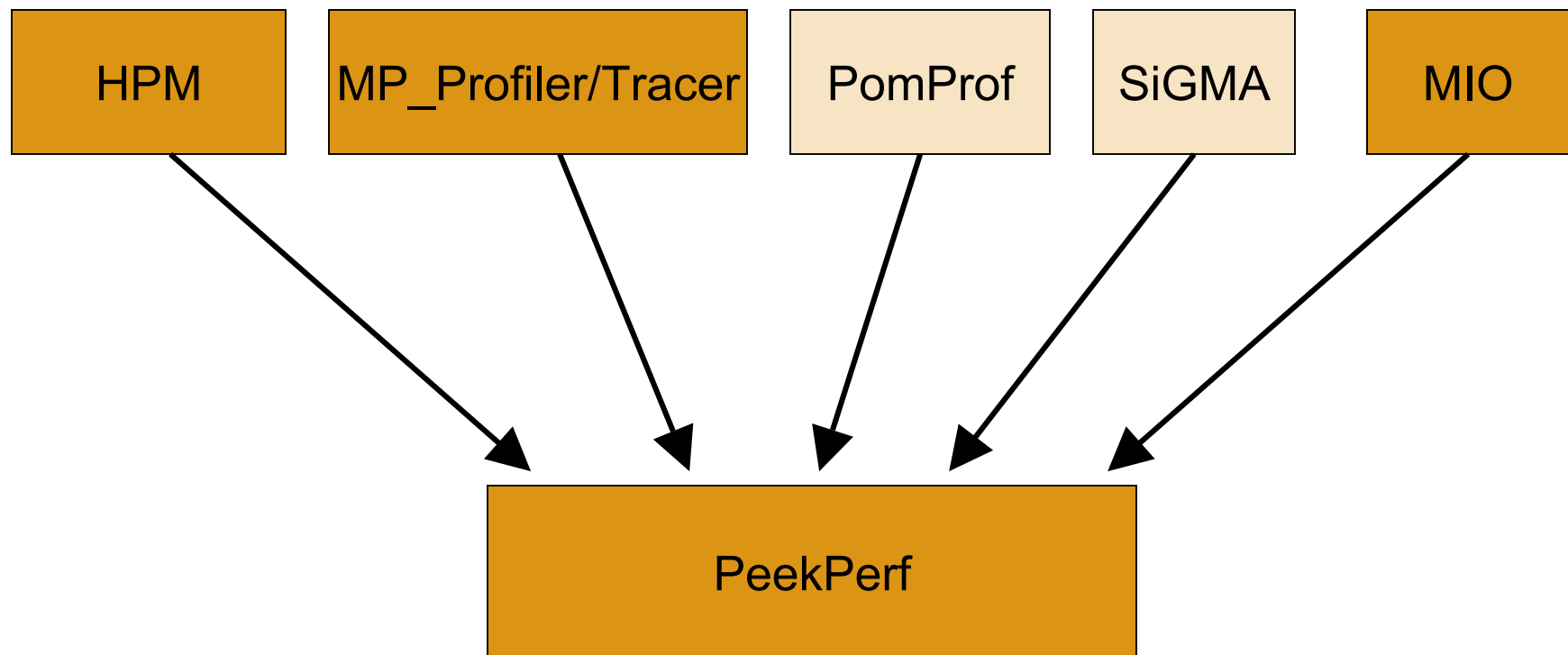I-Hsin Chung          ihchung@us.ibm.com
Guojing Cong          gcong@us.ibm.com
David Klepacki        klepacki@us.ibm.com

# Outline

- **IBM High Performance Computing Toolkit**

  – MPI performance: MP_Profiler

  – CPU performance: Xprofiler, HPM

  – Modular I/O: MIO

  – Visualization and analysis: PeekPerf

- **Related Tools**

- **Challenges**

# IBM HPCT Overview

| HPM | MP_Profiler/Tracer | PomProf | SiGMA | MIO |
|-----|--------------------|---------| ------|-----|

PeekPerf

# Message-Passing Performance:

- **MP_Profiler Library**

  – Captures "summary" data for MPI calls

  – Source code traceback

  – User **MUST** call MPI_Finalize() in order to get output files.

  – No changes to source code
    - MUST compile with –g to obtain source line number information

- **MP_Tracer Library**

  – Captures "timestamped" data for MPI calls

  – Source traceback

# MP_Profiler Summary Output

```
-----------------------------------------------------------
MPI Routine              #calls      avg. bytes      time(sec)
-----------------------------------------------------------
MPI_Comm_size                3           0.0          0.000
MPI_Comm_rank            12994           0.0          0.016
MPI_Send                19575       11166.9         13.490
MPI_Isend              910791        5804.2          9.216
MPI_Recv               138173        2767.9         73.835
MPI_Irecv              784936       15891.6          2.407
MPI_Sendrecv           894809         352.0         88.705
MPI_Wait              1537375           0.0        288.049
MPI_Waitall             44042           0.0         25.312
MPI_Bcast                 464       41936.8          3.272
MPI_Barrier              1312           0.0         34.206
MPI_Gather                 68       16399.1          2.680
MPI_Scatter                 6       17237.3          0.532
-----------------------------------------------------------
total communication time = 770.424 seconds.
total elapsed time       = 1168.662 seconds.
user cpu time            = 1160.960 seconds.
system time              = 0.620 seconds.
maximum memory size      = 68364 KBytes.

To check load balance :  grep "total comm" mpi_profile.*
```

# MP_Profiler Sample Call Graph Output

```
communication time = 143.940 sec, parent = gwwrloc
     MPI Routine              #calls          time(sec)
     MPI_Barrier               2311           143.734
     MPI_Gatherv               2311             0.206
communication time = 137.823 sec, parent = f2drecv
     MPI Routine              #calls          time(sec)
     MPI_Recv                 91959           137.823
communication time = 108.960 sec, parent = puttsf
     MPI Routine              #calls          time(sec)
     MPI_Barrier              23607           106.821
     MPI_Gatherv              23607             2.139
communication time = 94.435 sec, parent = fft_tri_recv
     MPI Routine              #calls          time(sec)
     MPI_Recv                  6378            94.435
communication time = 83.836 sec, parent = fft2drecv
     MPI Routine              #calls          time(sec)
     MPI_Recv                 93003            83.836
```

# MP_Profiler Message Size Distribution

```
msglen =              8 bytes,    elapsed time = 0.0029 msec
msglen =              8 bytes,    elapsed time = 0.0029 msec
msglen =             32 bytes,    elapsed time = 0.0030 msec
msglen =             64 bytes,    elapsed time = 0.0028 msec
msglen =             96 bytes,    elapsed time = 0.0028 msec
msglen =            160 bytes,    elapsed time = 0.0029 msec
msglen =            240 bytes,    elapsed time = 0.0030 msec
msglen =            320 bytes,    elapsed time = 0.0030 msec
msglen =            400 bytes,    elapsed time = 0.0030 msec
msglen =            480 bytes,    elapsed time = 0.0031 msec

                      .  .  .

                      .  .  .

                      .  .  .

msglen =         640000 bytes,    elapsed time = 0.3616 msec
msglen =         720000 bytes,    elapsed time = 0.4019 msec
msglen =         800000 bytes,    elapsed time = 0.4630 msec
msglen =        1000000 bytes,    elapsed time = 0.6618 msec
```

# Hardware Performance Monitor (HPM)

- **Provides comprehensive reports of events that are critical to performance on IBM systems.**

- **Gather critical hardware performance metrics, e.g.**

  - Number of misses on all cache levels

  - Number of floating point instructions executed

  - Number of instruction loads that cause TLB misses

- **Helps to identify and eliminate performance bottlenecks.**

# HPM

- **hpmcount**

  - Starts application and provides

    - Wall clock time
    - Hardware performance counter information
    - Resource utilization statistics

- **libhpm**

  - Library for program (including multi-thread) section instrumentation.

- **hpmstat**

  - Provides system wide reports for root

# Hpmcount Example

bash-2.05a$ hpm_count swim

…… // program output

hpmcount (V 2.5.4) summary

| | |
|---|---|
| **Execution time (wall clock time)** | **: 7.378159 seconds** |
| **######## Resource Usage Statistics ########** | |
| **Total amount of time in user mode** | **: 0.010000 seconds** |
| **Average shared memory use in text segment** | **: 672 Kbytes*sec** |
| **………** | |
| **PM_FPU_FDIV (FPU executed FDIV instruction)** | **: 0** |
| **PM_FPU_FMA (FPU executed multiply-add instruction)** | **: 0** |
| **PM_CYC (Processor cycles)** | **: 54331072** |
| **PM_FPU_STF (FPU executed store instruction)** | **: 2172** |
| **PM_INST_CMPL (Instructions completed)** | **: 17928229** |
| **Utilization rate** | **: 0.446 %** |
| **Total load and store operations** | **: 0.004 M** |
| **MIPS** | **: 2.140** |
| **Instructions per cycle** | **: 0.330** |

# Libhpm Example

**call f_hpmstart( 30, "Loop 300" )**

C$OMP PARALLELDO

C$OMP&SHARED (ALPHA,M,N,U,V,P,UNEW,VNEW,PNEW,UOLD,VOLD,POLD)

C$OMP&SHARED (JS,JE)

C$OMP&PRIVATE (I,J)

```
   DO 300 J=js,je
   DO 300 I=1,M
   UOLD(I,J) = U(I,J)+ALPHA*(UNEW(I,J)-2.*U(I,J)+UOLD(I,J))
    ……
   P(I,J) = PNEW(I,J)
300 CONTINUE
```

**call f_hpmstop( 30 )**

# Modular I/O (MIO)

- **Addresses the need of application-level optimization for I/O.**

- **Analyze and tune I/O at the application level**
  - For example, when an application exhibits the I/O pattern of sequential reading of large files
  - MIO
    - Detects the behavior
    - Invokes its asynchronous prefetching module to prefetch user data.

- **Planned Integration into HPC Toolkit with PeekPerf capabilities**
  - Source code traceback
  - Future capability for dynamic I/O instrumentation

# MP_Profiler Visualization Using PeekPerf

# MP_Tracer Visualization Using PeekPerf

# HPM Visualization Using PeekPerf

# Check Performance Counter
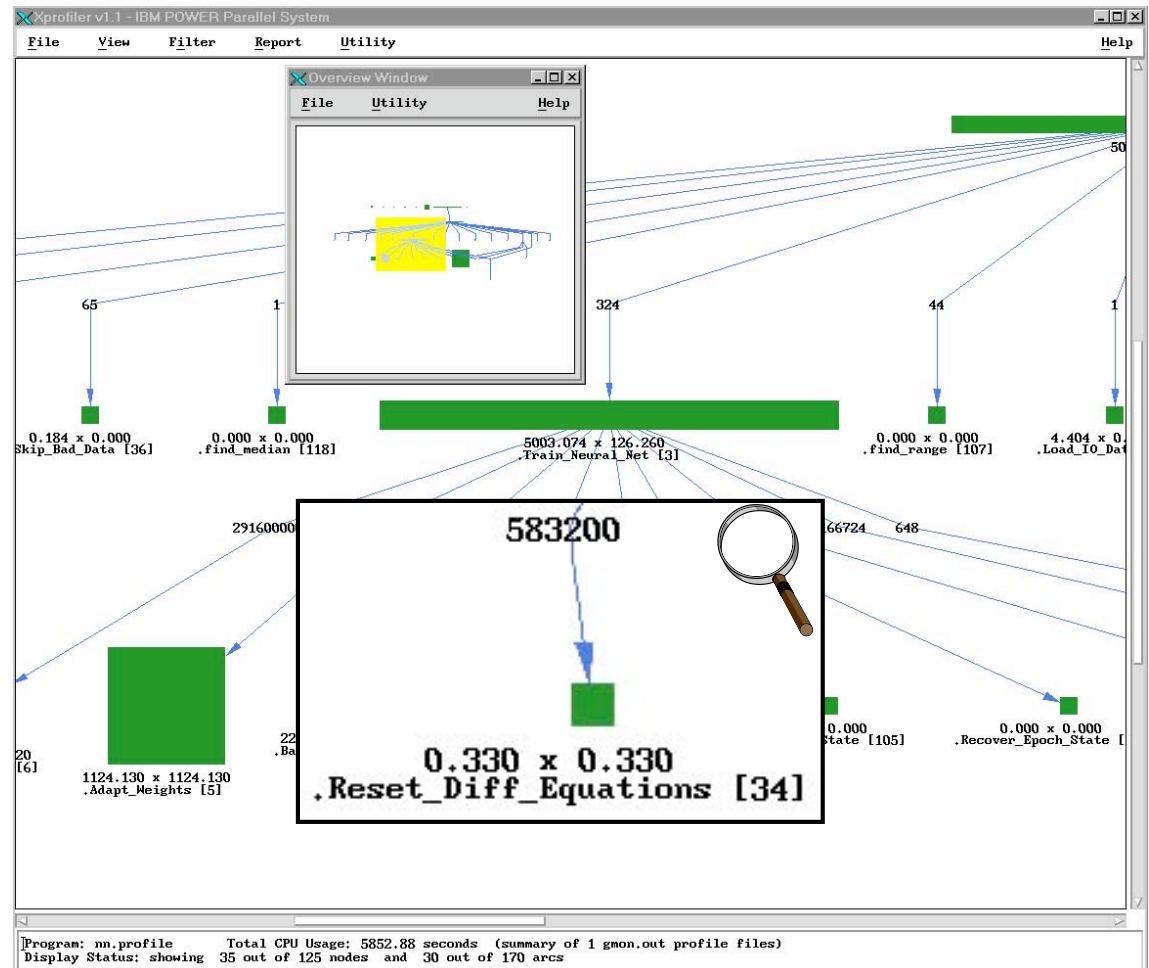
# Xprofiler

- **CPU profiling tool, extension of gprof**
  - can be used to profile both serial and parallel applications.

- **Graphical display of the call graphs of the application**

- **Provides quick access to the profiled data**

- **Helps users identify the CPU-intensive functions**

# Running Xprofiler

- **Compile the program with –pg**

- **Run the program**

- **gmon.out file is generated  (MPI applications generate gmon.out.1, …, gmon.out.n)**

- **Run Xprofiler**

# Xprofiler: Main Display

- **Width of a bar:** time including called routines

- **Height of a bar:** time excluding called routines

- **Call arrows labeled with number of calls**

- **Overview window for easy navigation (View → Overview)**

# Xprofiler: Flat Profile

- **Menu** `Report` **provides usual gprof reports plus some extra ones**

  - Flat Profile

  - Call Graph Profile

  - Function Index

  - Function Call Summary

  - Library Statistics

# Xprofiler: Source Code Window

- Source code Window displays source code with time profile (resolution 1 tick = 0.01 sec)

- Access
  - Select function in main display
  - → context menu
  - Select function in flat profile
  - → Code Display
  - → Show Source Code

# Xprofiler - Application View

# Related Work

- **PARAVER (UPC)**

  – Performance Visualization and Analysis Tool

  – Flexibility to represent traces from different environment

  – MPI trace, HW performance counter info (based on PAPI)

- **PAPI (Univ. of Tennessee)**

  – **P**erformance **A**pplication **P**rogramming **I**nterface

  – Design, standardize, and implement a portable and efficient API to access the hardware performance counters

# Related Work (continued)

- **TAU (Univ. of Oregon)**

  – **T**uning and **A**nalysis **U**tilities

  – Program and performance analysis tool framework for high-performance parallel and distributed computing

- **mpiP (LLNL/ORNL)**

  – Lightweight profiling library for MPI applications

  – Only collects statistical information about MPI functions

  – With less overhead and have much less data than tracing tools

  – Only uses communication during report generation stage

# Challenges

- **Hardware**

  – Different performance counter set

- **Scalability**

  – 64k nodes

  – Tracking communications for each pair:

    - $(65,536)^2$ x48 bytes = 200 GB

  – Number of processes and events vs. number of screen pixels

# Scalability

- **Abstraction**

  – Aggregated performance data

  – Sampling

  – Easy to manage

  – May lose fidelity

- **Performance data organization**

  – Database with one time processing

  – Query, data mining, clustering… etc.

  – Precise information

  – Hard to analyze

# Summary

- **IBM High Performance Computing Toolkit**
  - Ported
    - PeekPerf, MP_Profiler
  - Work in progress
    - Xprofiler, HPM
  - Next target
    - Modular I/O: MIO
- **Related Tools**
- **Challenges**
  - Architecture, Scalability

# DPOMP

- **Dynamically instruments OpenMP applications**

- **Reports various performance related information**

  – timings, overhead of OpenMP constructs.

- **Modify binaries with performance instrumentation**

  – without requiring access to source codes or recompilation.

- **POMP (Performance monitoring interface for OpenMP) implementation based on dynamic probes**

# Simulation Guided Memory Analyzer (SiGMA)

- **Helps to understand precise memory references causing poor memory utilization**

- **Provides fine-grained information useful for**
  - Tuning loop kernels, understanding the cache behavior

- **Consists of**
  - A pre-execution tool that instruments all instructions that refer to memory locations,

  - A runtime data collection tool that performs compression of the stream of memory addresses generated by the instrumentation,

  - Analysis tools that process memory reference trace to provide programmers with tuning information.